

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

**Annotating Programs for Automatic Summary
Generation**

Inventor(s):

Yong Rui

Anoop Gupta

Alejandro Acero

ATTORNEY'S DOCKET NO. MS1-416US

1 **RELATED APPLICATIONS**

2 This application claims the benefit of U.S. Provisional Application No.
3 60/153,730, filed September 13, 1999, entitled "MPEG-7 Enhanced Multimedia
4 Access" to Yong Rui, Jonathan Grudin, Anoop Gupta, and Liwei He, which is
5 hereby incorporated by reference.

6 **TECHNICAL FIELD**

7 This invention relates to audio/video programming and rendering thereof,
8 and more particularly to annotating programs for automatic summary generation.

9 **BACKGROUND OF THE INVENTION**

10 Watching television has become a common activity for many people,
11 allowing people to receive important information (e.g., news broadcasts, weather
12 forecasts, etc.) as well as simply be entertained. While the quality of televisions
13 on which programs are rendered has improved, so too have a wide variety of
14 devices been developed and made commercially available that further enhance the
15 television viewing experience. Examples of such devices include Internet
16 appliances that allow viewers to "surf" the Internet while watching a television
17 program, recording devices (either analog or digital) that allow a program to be
18 recorded and viewed at a later time, etc.

19 Despite these advances and various devices, mechanisms for watching
20 television programs are still limited to two general categories: (1) watching the
21 program "live" as it is broadcast, or (2) recording the program for later viewing.
22 Each of these mechanisms, however, limits viewers to watching their programs in
23 the same manner as they were was broadcast (although possibly time-delayed).

Often times, however, people do not have sufficient time to watch the entirety of a recorded television program. By way of example, a sporting event such as a baseball game may take 2 or 2½ hours, but a viewer may only have ½ hour that he or she can spend watching the recorded game. Currently, the only way for the viewer to watch such a game is for the viewer to randomly select portions of the game to watch (e.g., using fast forward and/or rewind buttons), or alternatively use a "fast forward" option to play the video portion of the recorded game back at a higher speed than that at which it was recorded (although no audio can be heard). Such solutions, however, have significant drawbacks because it is extremely difficult for the viewer to know or identify which portions of the game are the most important for him or her to watch. For example, the baseball game may have only a handful of portions that are exciting, with the rest being uninteresting and not exciting.

The invention described below addresses these disadvantages, providing for annotating of programs for automatic summary generation.

SUMMARY OF THE INVENTION

Annotating programs for automatic summary generation is described herein.

In accordance with one aspect, audio/video programming content is made available to a receiver from a content provider, and meta data is made available to the receiver from a meta data provider. The content provider and meta data provider may be the same or different devices. The meta data corresponds to the programming content, and identifies, for each of multiple portions of the programming content, an indicator of a likelihood that the portion is an exciting

portion of the content. The meta data can be used, for example, to allow summaries of the programming content to be generated by selecting the portions having the highest likelihoods of being exciting portions.

According to another aspect, exciting portions of a sporting event are automatically identified based on sports-specific events and sports-generic events. The audio data of the sporting event is analyzed to identify sports-specific events (such as baseball hits if the sporting event is a baseball program) as well as sports-generic events (such as excited speech from an announcer). These sports-specific and sports-generic events are used together to identify the exciting portions of the sporting event.

According to another aspect, exciting segments of a baseball program are automatically identified. Various features are extracted from the audio data of the baseball program and selected features are input to an excited speech classification subsystem and a baseball hit detection subsystem. The excited speech classification subsystem identifies probabilities that segments of the audio data contain excited speech (e.g., from an announcer). The baseball hit detection subsystem identifies probabilities that multiple-frame groupings of the audio data include baseball hits. These two sets of probabilities are input to a probabilistic fusion subsystem that determines, based on both probabilities, a likelihood that each of the segments is an exciting portion of the baseball program. These probabilities can then be used, for example, to generate a summary of the baseball program.

1 **BRIEF DESCRIPTION OF THE DRAWINGS**

2 The present invention is illustrated by way of example and not limitation in
3 the figures of the accompanying drawings. The same numbers are used
4 throughout the figures to reference like components and/or features.

5 Fig. 1 shows a programming distribution and viewing system in accordance
6 with one embodiment of the invention;

7 Fig. 2 illustrates an example of a suitable operating environment in which
8 the invention may be implemented;

9 Fig. 3 illustrates an exemplary programming content delivery architecture
10 in accordance with certain embodiments of the invention;

11 Fig. 4 illustrates an exemplary automatic summary generation process in
12 accordance with certain embodiments of the invention;

13 Fig. 5 illustrates part of an exemplary audio clip and portions from which
14 features are extracted;

15 Fig. 6 illustrates exemplary baseball hit templates that may be used in
16 accordance with certain embodiments of the invention; and

17 Fig. 7 is a flowchart illustrating an exemplary process for rendering a
18 program summary to a user in accordance with certain embodiments of the
19 invention.

20 **DETAILED DESCRIPTION**

21 **General System**

22 Fig. 1 shows a programming distribution and viewing system 100 in
23 accordance with one embodiment of the invention. System 100 includes a video
24 and audio rendering system 102 having a display device including a viewing area

1 104. Video and audio rendering system 102 represents any of a wide variety of
2 devices for playing video and audio content, such as a traditional television
3 receiver, a personal computer, etc. Receiver 106 is connected to receive and
4 render content from multiple different programming sources. Although illustrated
5 as separate components, rendering system 102 may be combined with receiver 106
6 into a single component (e.g., a personal computer or television). Receiver 106
7 may also be capable of storing content locally, in either analog or digital format
8 (e.g., on magnetic tapes, a hard disk drive, optical disks, etc.).

9 While audio and video have traditionally been transmitted using analog
10 formats over the airwaves, current and proposed technology allows multimedia
11 content transmission over a wider range of network types, including digital
12 formats over the airwaves, different types of cable and satellite systems
13 (employing both analog and digital transmission formats), wired or wireless
14 networks such as the Internet, etc.

15 Fig. 1 shows several different physical sources of programming, including a
16 terrestrial television broadcasting system 108 which can broadcast analog or
17 digital signals that are received by antenna 110; a satellite broadcasting system 112
18 which can transmit analog or digital signals that are received by satellite dish 114;
19 a cable signal transmitter 116 which can transmit analog or digital signals that are
20 received via cable 118; and an Internet provider 120 which can transmit digital
21 signals that are received by modem 122 via the Internet (and/or other network)
22 124. Both analog and digital signals can include programming made up of audio,
23 video, and/or other data. Additionally, a program may have different components
24 received from different programming sources, such as audio and video data from
25 cable transmitter 116 but data from Internet provider 120. Other programming

1 sources might be used in different situations, including interactive television
2 systems.

3 As described in more detail below, programming content made available to
4 system 102 includes audio and video programs as well as meta data corresponding
5 to the programs. The meta data is used to identify portions of the program that are
6 believed to be exciting portions, as well as how exciting these portions are
7 believed to be relative to one another. The meta data can be used to generate
8 summaries for the programs, allowing the user to view only the portions of the
9 program that are determined to be the most exciting.

10

11 **Exemplary Operating Environment**

12 Fig. 2 illustrates an example of a suitable operating environment in which
13 the invention may be implemented. The illustrated operating environment is only
14 one example of a suitable operating environment and is not intended to suggest
15 any limitation as to the scope of use or functionality of the invention. Other well
16 known computing systems, environments, and/or configurations that may be
17 suitable for use with the invention include, but are not limited to, personal
18 computers, server computers, hand-held or laptop devices, multiprocessor systems,
19 microprocessor-based systems, programmable consumer electronics (e.g., digital
20 video recorders), gaming consoles, cellular telephones, network PCs,
21 minicomputers, mainframe computers, distributed computing environments that
22 include any of the above systems or devices, and the like.

23 Alternatively, the invention may be implemented in hardware or a
24 combination of hardware, software, and/or firmware. For example, one or more
25

1 application specific integrated circuits (ASICs) could be designed or programmed
2 to carry out the invention.

3 Fig. 2 shows a general example of a computer 142 that can be used in
4 accordance with the invention. Computer 142 is shown as an example of a
5 computer that can perform the functions of receiver 106 of Fig. 1, or of one of the
6 programming sources of Fig. 1 (e.g., Internet provider 120). Computer 142
7 includes one or more processors or processing units 144, a system memory 146,
8 and a bus 148 that couples various system components including the system
9 memory 146 to processors 144.

10 The bus 148 represents one or more of any of several types of bus
11 structures, including a memory bus or memory controller, a peripheral bus, an
12 accelerated graphics port, and a processor or local bus using any of a variety of
13 bus architectures. The system memory 146 includes read only memory (ROM)
14 150 and random access memory (RAM) 152. A basic input/output system (BIOS)
15 154, containing the basic routines that help to transfer information between
16 elements within computer 142, such as during start-up, is stored in ROM 150.
17 Computer 142 further includes a hard disk drive 156 for reading from and writing
18 to a hard disk, not shown, connected to bus 148 via a hard disk drive interface 157
19 (e.g., a SCSI, ATA, or other type of interface); a magnetic disk drive 158 for
20 reading from and writing to a removable magnetic disk 160, connected to bus 148
21 via a magnetic disk drive interface 161; and an optical disk drive 162 for reading
22 from and/or writing to a removable optical disk 164 such as a CD ROM, DVD, or
23 other optical media, connected to bus 148 via an optical drive interface 165. The
24 drives and their associated computer-readable media provide nonvolatile storage
25 of computer readable instructions, data structures, program modules and other data

1 for computer 142. Although the exemplary environment described herein employs
2 a hard disk, a removable magnetic disk 160 and a removable optical disk 164, it
3 will be appreciated by those skilled in the art that other types of computer readable
4 media which can store data that is accessible by a computer, such as magnetic
5 cassettes, flash memory cards, random access memories (RAMs), read only
6 memories (ROM), and the like, may also be used in the exemplary operating
7 environment.

8 A number of program modules may be stored on the hard disk, magnetic
9 disk 160, optical disk 164, ROM 150, or RAM 152, including an operating system
10 170, one or more application programs 172, other program modules 174, and
11 program data 176. A user may enter commands and information into computer
12 142 through input devices such as keyboard 178 and pointing device 180. Other
13 input devices (not shown) may include a microphone, joystick, game pad, satellite
14 dish, scanner, or the like. These and other input devices are connected to the
15 processing unit 144 through an interface 168 that is coupled to the system bus
16 (e.g., a serial port interface, a parallel port interface, a universal serial bus (USB)
17 interface, etc.). A monitor 184 or other type of display device is also connected to
18 the system bus 148 via an interface, such as a video adapter 186. In addition to the
19 monitor, personal computers typically include other peripheral output devices (not
20 shown) such as speakers and printers.

21 Computer 142 operates in a networked environment using logical
22 connections to one or more remote computers, such as a remote computer 188.
23 The remote computer 188 may be another personal computer, a server, a router, a
24 network PC, a peer device or other common network node, and typically includes
25 many or all of the elements described above relative to computer 142, although

only a memory storage device 190 has been illustrated in Fig. 2. The logical connections depicted in Fig. 2 include a local area network (LAN) 192 and a wide area network (WAN) 194. Such networking environments are commonplace in offices, enterprise-wide computer networks, intranets, and the Internet. In certain embodiments of the invention, computer 142 executes an Internet Web browser program (which may optionally be integrated into the operating system 170) such as the "Internet Explorer" Web browser manufactured and distributed by Microsoft Corporation of Redmond, Washington.

When used in a LAN networking environment, computer 142 is connected to the local network 192 through a network interface or adapter 196. When used in a WAN networking environment, computer 142 typically includes a modem 198 or other means for establishing communications over the wide area network 194, such as the Internet. The modem 198, which may be internal or external, is connected to the system bus 148 via a serial port interface 168. In a networked environment, program modules depicted relative to the personal computer 142, or portions thereof, may be stored in the remote memory storage device. It will be appreciated that the network connections shown are exemplary and other means of establishing a communications link between the computers may be used.

Computer 142 also includes a broadcast tuner 200. Broadcast tuner 200 receives broadcast signals either directly (e.g., analog or digital cable transmissions fed directly into tuner 200) or via a reception device (e.g., via antenna 110 or satellite dish 114 of Fig. 1).

Computer 142 typically includes at least some form of computer readable media. Computer readable media can be any available media that can be accessed by computer 142. By way of example, and not limitation, computer readable

1 media may comprise computer storage media and communication media.
2 Computer storage media includes volatile and nonvolatile, removable and non-
3 removable media implemented in any method or technology for storage of
4 information such as computer readable instructions, data structures, program
5 modules or other data. Computer storage media includes, but is not limited to,
6 RAM, ROM, EEPROM, flash memory or other memory technology, CD-ROM,
7 digital versatile disks (DVD) or other optical storage, magnetic cassettes, magnetic
8 tape, magnetic disk storage or other magnetic storage devices, or any other media
9 which can be used to store the desired information and which can be accessed by
10 computer 142. Communication media typically embodies computer readable
11 instructions, data structures, program modules or other data in a modulated data
12 signal such as a carrier wave or other transport mechanism and includes any
13 information delivery media. The term "modulated data signal" means a signal that
14 has one or more of its characteristics set or changed in such a manner as to encode
15 information in the signal. By way of example, and not limitation, communication
16 media includes wired media such as wired network or direct-wired connection,
17 and wireless media such as acoustic, RF, infrared and other wireless media.
18 Combinations of any of the above should also be included within the scope of
19 computer readable media.

20 The invention has been described in part in the general context of
21 computer-executable instructions, such as program modules, executed by one or
22 more computers or other devices. Generally, program modules include routines,
23 programs, objects, components, data structures, etc. that perform particular tasks
24 or implement particular abstract data types. Typically the functionality of the
25

1 program modules may be combined or distributed as desired in various
2 embodiments.

3 For purposes of illustration, programs and other executable program
4 components such as the operating system are illustrated herein as discrete blocks,
5 although it is recognized that such programs and components reside at various
6 times in different storage components of the computer, and are executed by the
7 data processor(s) of the computer.

8

9 **Content Delivery Architecture**

10 Fig. 3 illustrates an exemplary programming content delivery architecture
11 in accordance with certain embodiments of the invention. A client 220 receives
12 programming content including both audio/video data 222 and meta data 224 that
13 corresponds to the audio/video data 222. In the illustrated example, an
14 audio/video data provider 226 is the source of audio/video data 222 and a meta
15 data provider 228 is the source of meta data 224. Alternatively, meta data 224 and
16 audio/video data 222 may be provided by the same source, or alternatively three or
17 more different sources.

18 The data 222 and 224 can be made available by providers 226 and 228 in
19 any of a wide variety of formats. In one implementation, data 222 and 224 are
20 formatted in accordance with the MPEG-7 (Moving Pictures Expert Group)
21 format. The MPEG-7 format standardizes a set of Descriptors (Ds) that can be
22 used to describe various types of multimedia content, as well as a set of
23 Description Schemes (DSs) to specify the structure of the Ds and their
24 relationship. In MPEG-7, the audio and video data 222 are each described as one
25 or more Descriptors, and the meta data 224 is described as a Description Scheme.

1 Client 220 includes one or more processor(s) 230 and renderer(s) 232.
2 Processor 230 receives audio/video data 222 and meta data 224 and performs any
3 necessary processing on the data prior to providing the data to renderer(s) 232.
4 Each renderer 232 renders the data it receives in a human-perceptive manner (e.g.,
5 playing audio data, displaying video data, etc.). The processing of data 222 and
6 224 can vary, and can include, for example, separating the data for delivery to
7 different renderers (e.g., audio data to a speaker and video data to a display
8 device), determining which portions of the program are most exciting based on the
9 meta data (e.g., probabilities included as the meta data), selecting the most
10 exciting segments based on a user-desired summary presentation time (e.g., the
11 user wants a 20-minute summary), etc.

12 Client 220 is illustrated as separate from providers 226 and 228. This
13 separation can be small (e.g., across a LAN) or large (e.g., a remote server located
14 in another city or state). Alternatively, data 222 and/or 224 may be stored locally
15 by client 220, either on another device such as an analog or digital video recorder
16 (not shown) coupled to client 220 or within client 220 (e.g., on a hard disk drive).

17 A wide variety of meta data 224 can be associated with a program. In the
18 discussions below, meta data 224 is described as being "excited segment
19 probabilities" which identify particular segments of the program and a
20 corresponding probability or likelihood that each segment is an "exciting"
21 segment. An exciting segment is a segment of the program believed to be
22 typically considered exciting to viewers. By way of example, baseball hits are
23 believed to be typically considered exciting segments of a baseball program.

24 The excited segment probabilities in meta data 224 can be generated in any
25 of a variety of manners. In one implementation, the excited segment probabilities

1 are generated manually (e.g., by a producer or other individual(s) watching the
2 program and identifying the exciting segments and assigning the corresponding
3 probabilities). In another implementation, the excited segment probabilities are
4 generated automatically by a process described in more detail below.
5 Additionally, the excited segment probabilities can be generated after the fact
6 (e.g., after a baseball game is over and its entirety is available on a recording
7 medium), or alternatively on the fly (e.g., a baseball game may be monitored and
8 probabilities generated as the game is played).

9

10 **Automatic Summary Generation**

11 The automatic summary generation process described below refers to
12 sports-generic and sports-specific events, and refers specifically to the example of
13 a baseball program. Alternatively, summaries can be automatically generated in
14 an analogous manner for other programs, including other sporting events.

15 The automatic summary generation process analyzes the audio data of the
16 baseball program and attempts to identify segments that include speech, and of
17 those segments which can be identified as being "excited" speech (e.g., the
18 excitement in an announcer's voice). Additionally, based on the audio data
19 segments that include baseball hits are also identified. These excited speech
20 segments and baseball hit segments are then used to determine, for each of the
21 excited speech segments, a probability that the segment is truly an exciting
22 segment of the program. Given these probabilities, a summary of the program can
23 be generated.

24 Fig. 4 illustrates an exemplary automatic summary generation process in
25 accordance with certain embodiments of the invention. The generation process

1 begins with the raw audio data 250 (also referred to as a raw audio clip), such as
2 the audio portion of data 222 of Fig. 3. The raw audio data 250 is the audio
3 portion of the program for which the summary is being automatically generated.
4 The audio data 250 is input to feature extractor 252 which extracts various features
5 from portions of audio data 250. In one implementation, feature extractor 252
6 extracts one or more of energy features, phoneme-level features, information
7 complexity features, and prosodic features.

8 Fig. 5 illustrates part of an exemplary audio clip and portions from which
9 features are extracted. Audio clip 258 is illustrated. Audio features are extracted
10 from audio clip 258 using two different resolutions: a sports-specific event
11 detection resolution used to assist in the identification of potentially exciting
12 sports-specific events, and a sports-generic event detection resolution used to
13 assist in the identification of potentially exciting sports-generic events. In the
14 illustrated example, the sports-specific event detection resolution is 10
15 milliseconds (ms), while the sports-generic event detection resolution is 0.5
16 seconds. Alternatively, other resolutions could be used.

17 As used herein, the sports-specific event detection is based on 10 ms
18 "frames", while the sports-generic event detection is based on 0.5 second
19 "windows". As illustrated in Fig. 5, the 10 ms frames are non-overlapping and the
20 0.5 second windows are non-overlapping, although the frames overlap the
21 windows (and vice versa). Alternatively, the frames may overlap other frames,
22 and/or the windows may overlap other windows.

23 Returning to Fig. 4, feature extractor 252 extracts different features from
24 audio data 250 based on both frames and windows of audio data 250. Exemplary
25 features which can be extracted by feature extractor 252 are discussed below.

1 Different embodiments can use different combinations of these features, or
2 alternatively use only selected ones of the features or additional features.

3 Extractor 252 extracts energy features for each of the 10ms frames of audio
4 data 250, as well as for each of the 0.5 second windows. For each frame or
5 window, feature vectors having, for example, one element are extracted that
6 identify the short-time energy in each of multiple different frequency bands. The
7 short-time energy for each frequency band is the average waveform amplitude in
8 the frequency band over the given time period (e.g., 10ms frame or 0.5 second
9 window). In one implementation, four different frequency bands are used: 0hz –
10 630hz, 630hz – 1720hz, 1720hz – 4400hz, and 4400hz and above, referred to as
11 E_1 , E_2 , E_3 , and E_4 , respectively. An additional feature vector is also calculated as
12 the summation of E_2 and E_3 , referred to as E_{23} .

13 The energy features extracted for each of the 10ms frames are also used to
14 determine energy statistics regarding each of the 0.5 second windows. Exemplary
15 energy statistics extracted for each frequency band E_1 , E_2 , E_3 , E_4 , and E_{23} for the
16 0.5 second window are illustrated in Table I.

17 Table I

18 Statistic	19 Description
20 maximum energy	21 The highest energy value of the frames 22 in the window.
23 average energy	24 The average energy value of the frames 25 in the window.
26 energy dynamic range	27 The energy range over the frames in the 28 window (the difference between the 29 maximum energy value and a minimum 30 energy value).

1 Extractor 252 extracts phoneme-level features for each of the 10ms frames
2 of audio data 250. For each frame, two well-known feature vectors are extracted:
3 a Mel-frequency Cepstral coefficient (MFCC) and the first derivative of the
4 MFCC (referred to as the delta MFCC). The MFCC is the *cosine* transform of the
5 pitch of the frame on the "Mel-scale", which is a gradually warped linear spectrum
6 (with coarser resolution at high frequencies).

7 Extractor 252 extracts information complexity features for each of the 10
8 ms frames of audio data 250. For each frame, a feature vector representing the
9 entropy (*Etr*) of the frame is extracted. For an *N*-point Fast Fourier Transform
10 (FFT) of an audio signal $s(t)$, with $S(n)$ representing the *n*th frequency's
11 component, entropy is defined as:

$$12 \quad Etr = \sum_{n=1}^N P_n \log P_n$$

13 where:

$$14 \quad P_n = \frac{|S(n)|^2}{\sum_{n=1}^N |S(n)|^2}$$

15 Extracting feature vectors representing entropy is well-known to those
16 skilled in the art and thus will not be discussed further except as it relates to the
17 present invention.

18 Extractor 252 extracts prosodic features for each of the 0.5 second windows
19 of audio data 250. For each window, a feature vector representing the pitch (*Pch*)
20 of the window is extracted. A variety of different well-known approaches can be
21

1 used in determining pitch, such as the auto-regressive model, the average
2 magnitude difference function, the maximum *a posteriori* (MAP) approach, etc.

3 The pitch is also determined for each 10ms frame of the 0.5 second
4 window. These individual frame pitches are then used to extract pitch statistics
5 regarding the pitch of the window. Exemplary pitch statistics extracted for each
6 0.5 second window are illustrated in Table II.

7 Table II

8 Statistic	9 Description
9 non-zero pitch count	The number of frames in the window that have a non-zero pitch value.
10 maximum pitch	The highest pitch value of the frames in the window.
11 minimum pitch	The lowest pitch value of the frames in the window.
12 average pitch	The average pitch value of the frames in the window.
13 pitch dynamic range	The pitch range over the frames in the window (the difference between the maximum and minimum pitch values).

16
17 Selected ones of the extracted features are passed by feature extractor 252
18 to an excited speech classification subsystem 260 and a baseball hit detection
19 subsystem 262. Excited speech classification subsystem 260 attempts to identify
20 segments of the audio data that include excited speech (sports-generic events),
21 while baseball hit detection subsystem 262 attempts to identify segments of the
22 audio data that include baseball hits (sports-specific events). The segments
23 identified by subsystems 260 and 262 may be of the same or alternatively different
24 sizes (and may be varying sizes). Probabilities generated for the segments are then
25

1 input to a probabilistic fusion subsystem 264 to determine a probability that the
2 segments are exciting.

3 Excited speech classification subsystem 260 uses a two-stage process to
4 identify segments of excited speech. In a first stage, energy and phoneme-level
5 features 266 from feature extractor 252 are input to a speech detector 268 that
6 identifies windows of the audio data that include speech (speech windows 270).
7 In the illustrated example, speech detector 268 uses both the E_{23} and the delta
8 MFCC feature vectors. For each 0.5 second window, if the E_{23} and delta MFCC
9 vectors each exceed corresponding thresholds, the window is identified as a
10 speech window 270; otherwise, the window is classified as not including speech.
11 In one implementation, the thresholds used by speech detector 268 are 2.0 for the
12 delta MFCC feature, and $0.07 * E_{cap}$ for the E_{23} feature (where E_{cap} is the highest
13 E_{23} value of all the frames in the audio clip (or alternatively all of the frames in the
14 audio clip that have been analyzed so far), although different thresholds could
15 alternatively be used.

16 In alternative embodiments, speech detector 268 may use different features
17 to classify segments as speech or not speech. By way of example, energy only
18 may be used (e.g., the window is classified as speech only if E_{23} exceeds a
19 threshold amount (such as $0.2 * E_{cap}$). By way of another example, energy and
20 entropy features may both be used (e.g., the window is classified as speech only if
21 the product of E_{23} and E_{tr} exceeds a threshold amount (such as 50,000).

22 In the second stage, pitch and energy features 272, received from feature
23 extractor 252, for each of the speech windows 270 are used by excited speech
24 classifier 274 to determine a probability that each speech window 270 is excited
25 speech. Classifier 274 then combines these probabilities to identify a probability

1 that a group of these windows (referred to as a segment, which in one
2 implementation is five seconds) is excited speech. Classifier 274 outputs an
3 indication of these excited speech segments 276, along with their corresponding
4 probabilities, to probabilistic fusion subsystem 264.

5 Excited speech classifier 274 uses six statistics regarding the energy E_{23}
6 features and the pitch (Pch) features extracted from each speech window 270:
7 maximum energy, average energy, energy dynamic range, maximum pitch,
8 average pitch, and pitch dynamic range. Classifier 274 concatenates these six
9 statistics together to generate a feature vector (having nine elements or
10 dimensions) and compares the feature vector to a set of training vectors (based on
11 corresponding features of training sample data) in two different classes: an
12 excited speech class and a non-excited speech class. The *posterior* probability of a
13 feature vector X (for a window 270) being in a class C_i , where C_1 is the class of
14 excited speech and C_2 is the class of non-excited speech, can be represented as:
15 $P(C_i | X)$. The probability of error in classifying the feature vector X can be
16 reduced by classifying the data to the class having the *posterior* probability that is
17 the highest.

18 Speech classifier 274 determines the *posterior* probability $P(C_i | X)$ using
19 learning machines. A wide variety of different learning machines can be used to
20 determine the *posterior* probability $P(C_i | X)$. Three such learning machines are
21 described below, although other learning machines could alternatively be used.

22 The *posterior* probability $P(C_i | X)$ can be determined using parametric
23 machines, such as Bayes rule:

$$P(C_i | X) = \frac{P(C_i)p(X | C_i)}{p(X)}$$

1 where $p(X)$ is the data density, $P(C_i)$ is the prior probability, and $p(X | C_i)$ is the
2 conditional class density. The data density $p(x)$ is a constant for all the classes and
3 thus does not contribute to the decision rule. The prior probability $P(C_i)$ can be
4 estimated from labeled training data (e.g., excited speech and non-excited speech)
5 in a conventional manner. The conditional class density $p(X | C_i)$ can be
6 calculated in a variety of different manners, such as the Gaussian (Normal)
7 distribution $N(\mu, \sigma)$. The μ parameter (mean) and the σ parameter (standard
8 deviation) can be determined using the well-known Maximum Likelihood
9 Estimation (MLE):

$$10 \quad 11 \quad \mu = \frac{1}{n} \sum_{k=1}^n X_k$$

$$12 \quad 13 \quad \sigma^2 = \frac{1}{n} \sum_{k=1}^n (X_k - \mu)^2$$

14 where n is the number of training samples and X represents the training samples.
15

16 Another type of machines that can be used to determine the *posterior*
17 probability $P(C_i | X)$ are non-parametric machines. The K nearest neighbor
18 technique is an example of such a machine. Using the K nearest neighbor
19 technique:

$$20 \quad 21 \quad P(C_i | X) = \frac{\frac{K_i}{nV}}{\sum_i \frac{K_i}{nV}} = \frac{K_i}{K}$$

22 where V is the volume around feature vector X , V covers K labeled (training)
23 samples, and K_i is the number of samples in class C_i .
24

1 Another type of machines that can be used to determine the *posterior*
2 probability $P(C_i | X)$ are semi-parametric machines, which combine the advantages
3 of non-parametric and parametric machines. Examples of such semi-parametric
4 machines include Gaussian mixture models, neural networks, and support vector
5 machines (SVMs).

6 Any of a wide variety of well-known training methods can be used to train
7 the SVM. After the SVM is trained, a sigmoid function is trained to map the SVM
8 outputs into *posterior* probabilities. The *posterior* probability $P(C_i | X)$ can then
9 be determined as follows:

$$10 \quad P(C_i | X) = \frac{1}{1 + \exp(AX + B)}$$

11 where A and B are the parameters of the sigmoid function. The parameters A and
12 B are determined by reducing the negative log likelihood of training data (f_i, t_i) ,
13 which is a cross-entropy error function:
14

$$16 \quad \text{min} - \sum_i t_i \log(p_i) + (1 - t_i) \log(1 - p_i)$$

17 where
18

$$20 \quad p_i = \frac{1}{1 + \exp(Af_i + B)}$$

21 The cross-entropy error function minimization can be performed using any
22 number of conventional optimization processes. The training data (f_i, t_i) can be the
23 same training data used to train the SVM, or other data sets. For example, the
24

1 training data (f_i, t_i) can be a hold out set (in which a fraction of the initial training
2 set, such as 30%, is not used to train the SVM but is used to train the sigmoid) or
3 can be generated using three-fold cross-validation (in which the initial training set
4 is split into three parts, each of three SVMs is trained on permutations of two out
5 of three parts, and the f_i are evaluated on the remaining third, and the union of all
6 three sets f_i forming the training set of the sigmoid).

7 Additionally, an out-of-sample model is used to avoid "overfitting" the
8 sigmoid. Out-of-sample data is modeled with the same empirical density as the
9 sigmoid training data, but with a finite probability of opposite label. In other
10 words, when a positive example is observed at a value f_i , rather than using $t_i=1$, it
11 is assumed that there is a finite chance of opposite label at the same f_i in the out-
12 of-sample data. Therefore, a value of $t_i=1-\epsilon_+$ is used, for some ϵ_+ . Similarly, a
13 negative example will use a target value of $t_i=\epsilon_-$.

14 Regardless of the manner in which the *posterior* probability $P(C_i | X)$ for a
15 0.5 second window is determined, the *posterior* probabilities for multiple windows
16 are combined to determine the *posterior* probability for a segment. In one
17 implementation, each segment is five seconds, so the *posterior* probabilities of ten
18 adjacent windows are used to determine the *posterior* probability for each
19 segment.

20 The *posterior* probabilities for the multiple windows can be combined in a
21 variety of different manners. In one implementation, the *posterior* probability of
22 the segment being an exciting segment, referred to as $P(ES)$, is determined by
23 averaging the *posterior* probabilities of the windows in the segment:

$$24 P(ES) = \frac{1}{M} \sum_{m=1}^M P(C_1 | X_m)$$

25

1 where C_1 represents the excited speech class and M is the number of windows in
2 the segment.

3 Which ten adjacent windows to use for a segment can be determined in a
4 wide variety of different manners. In one implementation, if ten or more adjacent
5 windows include speech, then those adjacent windows are combined into a single
6 segment (e.g., which may be greater than ten windows, or, if too large, which may
7 be pared down into multiple smaller ten-window segments). However, if there are
8 fewer than ten adjacent windows, then additional windows are added (before
9 and/or after the adjacent windows, between multiple groups of adjacent windows,
10 etc.) to get the full ten windows, with the *posterior* probability for each of these
11 additional windows being zero.

12 The probabilities $P(ES)$ of these segments including excited speech 276 (as
13 well as an indication of where these segments occur in the raw audio clip 250) are
14 then made available to probabilistic fusion subsystem 264. Subsystem 264
15 combines the probabilities 276 with information received from baseball hit
16 detection subsystem 262, as discussed in more detail below.

17 Baseball hit detection subsystem 262 uses energy features 278 from feature
18 extractor 252 to identify baseball hits within the audio data 250. In one
19 implementation, the energy features 278 include the E_{23} and E_4 features discussed
20 above. Two additional features are also generated, which may be generated by
21 feature extractor 252 or alternatively another component (not shown). These
22 additional features are referred to as ER_{23} and ER_4 , and are discussed in more
23 detail below.

24 Hit detection is performed by subsystem 262 based on 25-frame groupings.
25 A sliding selection of 25 consecutive 10ms frames of the audio data 250 is

1 analyzed, with the frame selection sliding frame-by-frame through the audio data
2 250. The features of the 25-frame groupings and a set of hit templates 280 are
3 input to template matcher 282. Template matcher 282 compares the features of
4 each 25-frame grouping to the hit templates 280, and based on this comparison
5 determines a probability as to whether the particular 25-frame grouping contains a
6 hit. An identification of the 25-frame groupings (e.g., the first frame in the
7 grouping) and their corresponding probabilities are output by template matcher
8 282 as hit candidates 284.

9 Multiple-frame groupings are used to identify hits because the sound of a
10 baseball hit is typically longer in duration than a single frame (which is, for
11 example, only 10 ms). The baseball hit templates 280 are established to capture
12 the shape of the energy curves (using the four energy features discussed above)
13 over the time of the groupings (e.g., 25 10ms frames, or 0.25 seconds). Baseball
14 hit templates 280 are designed so that the hit peak (the energy peak) is at the 8th
15 frame of the 25-frame grouping. The additional features ER_{23} and ER_4 are
16 calculated by normalizing the E_{23} and E_4 features based on the energy features in
17 the 8th frame as follows:

$$ER_{23}(i) = \frac{E_{23}(i)}{E_{23}(8)}$$

$$ER_4(i) = \frac{E_4(i)}{E_4(8)}$$

22 where i ranges from 1 to 25, $E_{23}(8)$ is the E_{23} energy in the 8th frame, and $E_4(8)$ is
23 the E_4 energy in the 8th frame.
24

1 Fig. 6 illustrates exemplary baseball hit templates 280 that may be used in
2 accordance with certain embodiments of the invention. The templates 280 in
3 Fig. 6 illustrate the shape of the energy curves over time (25 frames) for each of
4 the four features E_{23} , E_4 , ER_{23} , and ER_4 .

5 For each group of frames, template matcher 282 determines the probability
6 that the group contains a baseball hit. This can be accomplished in multiple
7 different manners, such as un-directional or directional template mapping.
8 Initially, the four feature vectors for each of the 25 frames are concatenated,
9 resulting in a 100-element vector. The templates 280 are similarly concatenated
10 for each of the 25 frames, also resulting in a 100-element vector. The probability
11 of a baseball hit in a grouping $P(HT)$ can be calculated based on the Mahalanobis
12 distance D between the concatenated feature vector and the concatenated template
13 vector as follows:

$$14 \quad D^2 = (\vec{X} - \vec{T})^T \Sigma^{-1} (\vec{X} - \vec{T})$$
$$15$$

16 where \vec{X} is the concatenated feature vector, \vec{T} is the concatenated template vector,
17 and Σ is the covariance matrix of \vec{T} . Additionally, Σ is restricted to being a
18 diagonal matrix, allowing the baseball hit probability $P(HT)$ to be determined as
19 follows:

$$20$$
$$21 \quad P(HT) = \frac{\exp(-\frac{1}{2}D^2)}{C + \exp(-\frac{1}{2}D^2)}$$
$$22$$
$$23$$

1 where C is a constant that is data dependent (e.g., $\exp(-0.5D^2)$), where D^2 is the
2 distance between the concatenated feature vector and a template for non-hit
3 signals).

4 Alternatively, a directional template matching approach can be used, with
5 the distance D being calculated as follows:

6

$$7 D^2 = (\vec{X} - \vec{T})^T I \times \Sigma^{-1} (\vec{X} - \vec{T})$$

8 where I is a diagonal indicator matrix. The indicator matrix I is adjusted to
9 account for over-mismatches or under-mismatches (an over-mismatch is actually
10 good). In one implementation, when the values of E_{23} for the 25-frame grouping
11 are overmatching the templates (e.g., more than a certain number (such as one-
12 half) of the data values in the 25-frame grouping are higher than the corresponding
13 template values), then $I = \text{diag}[1, \dots, 1, -1, 1, \dots, 1]$ where the -1 is at location 8.
14 However, when the values of E_{23} for the 25-frame grouping are under-matching
15 the templates (e.g., less than a certain number (such as one-half) of the data values
16 in the 25-frame grouping are less than the corresponding template values), then $I =$
17 $\text{diag}[-1, \dots, -1, -1, -1, \dots, -1]$ where the 1 is at location 8.

18 Although hit detection is described as being performed across all of the
19 audio data 250, alternatively hit detection may be performed on only selected
20 portions of the audio data 250. By way of example, hit detection may only be
21 performed on the portions of audio data 250 that are excited speech segments (or
22 speech windows) and for a period of time (e.g., five seconds) prior to those excited
23 speech segments (or speech windows).

1 Probabilistic fusion generator 286 of subsystem 264 receives the excited
2 speech segment probabilities $P(ES)$ from excited speech classification subsystem
3 260 and the baseball hit probabilities $P(HT)$ from baseball hit detection subsystem
4 262 and combines those probabilities to identify probabilities $P(E)$ that segments
5 of the audio data 250 are exciting. Probabilistic fusion generator 286 searches for
6 hit frames within the 5-second interval of the excited speech segment. This
7 combining is also referred to herein as "fusion".

8 Two different types of fusion can be used: weighted fusion and conditional
9 fusion. Weighted fusion applies weights to each of the probabilities $P(ES)$ and
10 $P(HT)$ adds the results to obtain the value $P(E)$ as follows:

$$11 \quad P(E) = W_{ES}P(ES) + W_{HT}P(HT)$$

13 where the weights W_{ES} and W_{HT} sum up to 1.0. In one implementation, W_{ES} is 0.83
14 and W_{HT} is 0.17, although other weights could alternatively be used.

15 Conditional fusion, on the other hand, accounts for the detected baseball
16 hits adjusting the confidence level of the $P(ES)$ estimation (e.g., that the excited
17 speech probability is not high due to mislabeling a car horn as speech). The
18 conditional fusion is calculated as follows:

$$20 \quad P(E) = P(CF)P(ES)$$
$$21 \quad P(CF) = P(CF | HT)P(HT) + P(CF | \overline{HT})P(\overline{HT})$$
$$22 \quad P(\overline{HT}) = 1 - P(HT)$$

23 where $P(CF)$ is the probability of how much confidence there is in the $P(ES)$
24 estimation, and $P(\overline{HT})$ is the probability that there is no hit. $P(CF|HT)$ represents
25 the probability that we are confident that $P(ES)$ is accurate given there is a

baseball hit. Similarly, $P(CF|\bar{HT})$ represents the probability that we are confident that $P(ES)$ is accurate given there is no baseball hit. Both conditional probabilities $P(CF|HT)$ and $P(CF|\bar{HT})$ can be estimated from the training data. In one implementation, the value of $P(CF|HT)$ is 1.0 and the value of $P(CF|\bar{HT})$ is 0.3.

The final probability $P(E)$ that a segment is an exciting segment is then output by generator 286, identifying the exciting segments 288. These final probabilities, and an indication of the segments they correspond to, are stored as the meta data 224 of Fig. 3.

The actual portions of the program rendered for a user as the summary of the program are based on these exciting segments 288. Various modifications may be made, however, to make the rendering smoother. Examples of such modifications include: starting rendering of the exciting segment a period of time (e.g., three seconds) earlier than the hit (e.g., to render the pitching of the ball); merging together overlapping segments; merging together close-by (e.g., within ten seconds) segments; etc.

Once the probabilities that segments are exciting are identified, the user can choose to view a summary or highlights of the program. Which segments are to be delivered as the summary can be determined locally (e.g., on the user's client computer) or alternatively remotely (e.g., on a remote server).

Additionally, various "pre-generated" summaries may be generated and maintained by remote servers. For example, a remote server may identify which segments to deliver if a 15-minute summary is requested and which segments to deliver if a 30-minute summary is requested, and then store these identifications. By pre-generating such summaries, if a user requests a 15-minute summary, then

1 the pre-generated indications simply need to be accessed rather than determining,
2 at the time of request, which segments to include in the summary.

3 Fig. 7 is a flowchart illustrating an exemplary process for rendering a
4 program summary to a user in accordance with certain embodiments of the
5 invention. The acts of Fig. 7 may be implemented in software, and may be carried
6 out by a receiver 106 of Fig. 1 or alternatively a programming source of Fig. 1
7 (e.g., Internet provider 120).

8 Initially, the user request for a summary is received along with parameters
9 for the summary (act 300). The parameters of the summary identify what level of
10 summary the user desires, and can vary by implementation. By way of example, a
11 user may indicate as the summary parameters that he or she wants to be presented
12 with any segments that have a probability of 0.75 or higher of being exciting
13 segments. By way of another example, a user may indicate as the summary
14 parameters that he or she wants to be presented with a 20-minute summary of the
15 program.

16 The meta data corresponding to the program (the exciting segment
17 probabilities $P(E)$) is then accessed (act 302), and the appropriate exciting
18 segments identified based on the summary parameters (act 304). Once the
19 appropriate exciting segments are identified, they are rendered to the user (act
20 306). The manner in which the appropriate exciting segments are identified can
21 vary, in part based on the nature of the summary parameters. If the summary
22 parameters indicate that all segments with a $P(E)$ of 0.75 or higher should be
23 presented, then all segments with a $P(E)$ of 0.75 or greater are identified. If the
24 summary parameters indicate that a 20-minute summary should be generated, then
25 the appropriate segments are identified by determining (based on the $P(E)$ of the

1 segments and the lengths of the segments) the segments having the highest $P(E)$
2 that have a combined length less than 20 minutes.

3

4 **Conclusion**

5 Although the description above uses language that is specific to structural
6 features and/or methodological acts, it is to be understood that the invention
7 defined in the appended claims is not limited to the specific features or acts
8 described. Rather, the specific features and acts are disclosed as exemplary forms
9 of implementing the invention.